



**CERTIFIED COPY OF
PRIORITY DOCUMENT**

Patent Office
Canberra



I, GAYE TURNER, TEAM LEADER EXAMINATION SUPPORT AND SALES hereby certify that annexed is a true copy of the Provisional specification in connection with Application No. PR 0455 for a patent by DYNAMIC DIGITAL DEPTH RESEARCH PTY LTD filed on 29 September 2000.



WITNESS my hand this
Fourteenth day of August 2001

A handwritten signature in cursive script, appearing to read "G Turner".

GAYE TURNER
TEAM LEADER EXAMINATION
SUPPORT AND SALES

AUSTRALIA

Patents Act 1990

ORIGINAL

PROVISIONAL SPECIFICATION

IMAGE CONVERSION AND ENCODING TECHNIQUES

The invention is described in the following statement:

IMAGE CONVERSION AND ENCODING TECHNIQUES

Field of the Invention

The present invention is generally directed towards stereoscopic image synthesis and more particularly toward an improved method of converting two dimensional (2D) images for further encoding, transmission and decoding for the purpose of stereoscopic display.

Background of the Invention

A number of image processing tasks require that the depth of objects within an image be known. Such tasks include the application of special effects to film and video sequences and the conversion of 2D images into stereoscopic 3D. Determining the depth of objects may be referred to as the process of creating a depth map. In a depth map each object is coloured a shade of grey such that the shade indicates the depth of the object from a fixed point. Typically an object that is distant will be coloured in a dark shade of grey whilst a close object will be lighter. A standard convention for the creation of depth maps is yet to be adopted, and the reverse colouring may be used or different colours may be used to indicate different depths. For the purposes of explanation in this disclosure distant objects will be coloured darker than closer objects, and the colouring will typically be grey scale.

In order to obtain acceptable results the creation of a depth map from an existing 2D image has historically been undertaken manually. It will be appreciated that an image is merely a series of pixels to a computer, whereas a human operator is capable of distinguishing objects and their relative depths.

The creation of depth maps involves a system whereby each object of the image to be converted is outlined manually and a depth assigned to the object. This process is understandably slow, time consuming and costly. The outlining step is usually undertaken using a software program in conjunction with a mouse. Examples of a software program that may be used to undertake this task is Adobe "After Effects". An operator using Adobe would typically draw round the outline of each object that requires a depth to be assigned and then fill or "colour in" the object with the desired shades of grey that defines the depth or distance from the viewer required. This process would then be repeated for each object in

the image. Further, where a number of images are involved, for example a film, it will also be necessary to carry out these steps for each image or frame of the film.

In the traditional system the outline of the image would typically be described as some form of curve, for example a bezier curve. The use of such a curve enables the operator to alter the shape of the outline such that the outline of the object can be accurately aligned with the object.

Should a series of images require depth mapping e.g., a film or video, then the process would be repeated for each frame in the sequence.

It is likely that the size, position and/or depth of an object may change through a sequence. In this case the operator is required to manually track the object in each frame and processing each frame by correcting the curve, and updating the object depth by changing the shade of grey as necessary. It will be appreciated that this is a slow, tedious, time consuming and expensive process.

Previous attempts have been made to improve this process. The prior art describes techniques that attempt to automatically track the outline of the object as it moves from frame to frame. An example of such a technique is the application of Active Contours (ref: Active Contours – Andrew Blake and Michael Isard – ISBN 3-540-76217-5). The main limitation of this approach is the need to teach the software implementing the technique the expected motion of the object being tracked. This is a significant limitation when either the expected motion is not known, complex deformations are anticipated, or numerous objects with different motion characteristics are required to be tracked simultaneously.

Point-based tracking approaches have also been used to define the motion of outlines. These are popular in editing environments such as Commotion and After Effects. However, their application is very limited because it is frequently impossible to identify a suitable tracking point whose motion reflects the motion of the object as a whole. Point tracking is sometimes acceptable when objects are undergoing simple translations, but will not handle shape deformations, occlusions, or a variety of other common problems.

An Israeli company, AutoMedia, has produced a software product called AutoMasker. This enables an operator to draw the outline of an object and track it from frame to frame. The product relies on tracking the colour of an object and

thus fails when similar coloured objects intersect. The product also has difficulty tracking objects that change in size over subsequent frames, for example, as an object approaches a viewer or moves forward on the screen.

None of these approaches are able to acceptably assign, nor track, depth maps, and thus the creating of the depth maps is still a manual system.

Other techniques are described in the prior art and rely on reconstructing the movement of the camera originally used to record the 2D sequence. The limitation of these techniques is the need for camera motion within the original image sequence and the presence of well-defined features within each frame that can be used as tracking points.

Object of the Invention

Presently, it is necessary for an operator to manually create a depth map for each frame of an image, so as to obtain acceptable results. It is an object of the present invention to reduce the number of frames that require manual depth creation, thereby reducing the time commitments for operators creating the depth maps.

There remains a set of frames for which the depth maps are still to be created manually. It is a further object of the invention to assist the manual process of depth map creation for these frames.

Summary of the Invention

With the above objects in mind the present invention provides a method of creating a depth map including the steps of:

- assigning a depth to at least one pixel or portion of an image;
- determining x,y coordinates image characteristics for each said at least one pixel or portion of said image;
- utilising said depth(s), image characteristics and respective x,y coordinates to determine an algorithm to ascertain depth characteristics as a function of x,y coordinates and image characteristics;
- utilising said algorithm to calculate a depth characteristic for each pixel or portion of said image;
- wherein said depth characteristics form a depth map for said image.

In a further aspect the present invention provides a method of creating a series of depth maps for an image sequence including the steps of:

receiving at least one depth map for at least one frame of said image sequence;

utilising said at least one depth map to determine an algorithm to ascertain the depth characteristics as a function of x,y coordinates and image characteristics;

utilising said algorithm to create a depth map for each frame of said image sequence.

In yet a further aspect the present invention provides a method of creating a series of depth maps for an image sequence including the steps of:

selecting at least one key frame from said image sequence;

for each at least one key frame assigning a depth to at least one pixel or portion of each frame;

determining x,y coordinates and image characteristics for each said at least one pixel or portion of each said frame;

utilising said depth(s), image characteristics and respective x,y coordinates for each said at least one frame to determine an algorithm for each said at least one frame to ascertain depth characteristics as a function of x,y coordinates and depth characteristics;

utilising each algorithm to calculate depth characteristics for each pixel or portion of each said at least one frame;

wherein said depth characteristics form a depth map for each said at least one frame.

utilising each depth map to determine an algorithm to ascertain the depth characteristics for each frame as a function of x,y coordinates and image characteristics;

utilising said algorithm to create respective depth maps for each frame of said image sequence.

It will be understood that the system in referring to an algorithm may in fact create a number of different functions in order to create the depth maps as a result of the x,y coordinates and image characteristics.

A system implementing the present invention may elect to predetermine which frames in a sequence are to be considered key frames, for example each

fifth frame. The algorithm will also ideally consider time as an input to the algorithm to further refine the processing.

Brief Description of the Invention

The invention is intended to improve the process of producing depth maps for associated 2D images. This preferred embodiment involves the two phases of generating key-frame depth maps, and generating the remaining maps.

The first phase obtains a small amount of data from the user. This data is indicative of the basic structure of the scene. The 2D image and this associated data are presented to an algorithm that is capable of learning the relationship between the depth z assigned by the user to various image pixels, its x and y location, and image characteristics. The image characteristics include, although are not limited to, the RGB value for each pixel. In general the algorithm solves the equation

$$z = f(x, y, R, G, B)$$

for each pixel in the frame that the user has defined.

The algorithm then applies this learned relationship to the remaining pixels in the image to generate a depth map. If necessary, the user can refine their data to improve the accuracy of the depth map.

The second phase requires 2D images and associated depth maps to be provided at selected key-frames. The depth maps at these key-frames may be manually generated as previously disclosed by the applicants, or produced automatically using depth capture techniques including, although not limited to, laser range finders i.e. LIDAR (Light Direction And Range) devices and depth-from-focus techniques.

The 2D image and associated depth map(s), for each key-frame, is presented to an algorithm that is capable of learning the relationship between the depth z assigned to each pixel in the remaining frames, its x and y location and image characteristics. The image characteristics include, although are not limited to, the RGB value of each pixel. In general the algorithm solves the equation

$$z = f(x, y, R, G, B)$$

for each pixel in the key-frames.

The algorithm is then presented with each subsequent frame between the adjacent key-frames and for each pixel uses the algorithm to calculate the value of z .

In the Drawings

- 5 Figure 1 shows the training process of Phase One.
- Figure 2 shows the conversion process of Phase One.
- Figure 3 shows the training process of Phase Two.
- Figure 4 shows the conversion process of Phase Two.

Detailed Description of the Invention

- 10 The invention provides a technique for constructing 3D stereo information from one or more 2D images.

Phase One

- The first phase operates on a single image. A user is presented with the image and defines approximate depths for various regions in the image using a simple graphical interface. The graphical interface may provide various tools to assist the user in assigning depths to pixels, including but not limited to pen and paintbrush tools, area fill tools, tools that assign a depth based on the pixel colour, and tools that assign a depth based on the current pixel depth. The result of this process is that the depth is defined for a subset of the pixels in the image.

20 Create Mapping Function

- Once the system is provided with the image and some pixel depths, the system then observes and analyses the pixels with defined depths in order to create a mapping function. The mapping function may be a process or function that takes as input any measurement of a pixel or a set of pixels from the image and provides as output a depth value for the pixel or set of pixels.

- Individual pixel measurements may consist of red, green and blue values, or other measurements such as luminance, chrominance, contrast and spatial measurements such as horizontal and vertical positioning in the image. Alternatively the mapping function may operate on higher level image features, such as larger sets of pixels and measurements on a set of pixels such as mean and variance or edges, corners etc (i.e. the response of a feature detector). Larger sets of pixels may for example represent segments in the image, being sets of connected pixels forming a homogenous region.

For illustrative purposes only, a pixel may be represented in the form
 x, y, R, G, B, z

where x and y represent the x and y coordinates of the pixel, R, G, B represent the red, green and blue values of that pixel, and z represents the depth of that pixel.

5 Values of z are only defined where the user has specified a value.

The mapping function is learnt by capturing the relationship between image data and depth data for the pixels identified by the user. The mapping function may take the form of any generic processing unit, where input data is received, processed, and an output given. Preferably, this processing unit is amenable to
 10 a learning process, where its nature is determined by examination of the user data and corresponding image data.

The process of learning this relationship between input data, and desired output would be understood by those who have worked in the areas of artificial intelligence or machine learning, and may take on many forms. It is noted that
 15 these persons would not normally work in the areas of stereoscopic systems, or conversion of 2D images to 3D. In these fields, such mapping functions are known and include, although are not limited to, neural networks, decision trees, decision graphs, model trees and nearest-neighbour classifiers. Preferred embodiments of a learning algorithm are those that seek to design a mapping
 20 function that minimises some measurement of mapping error and that generalise satisfactorily for values outside the original data set.

The learning algorithm attempts to generalise the relationships between the 2D image information and the depth data specified by the user. This generalisation will then be applied to complete the depth maps for the entire
 25 sequence. Examples of successful learning algorithms are the back-propagation algorithm for learning neural networks, the C4.5 algorithm for learning decision trees, and the K-Means algorithm for learning cluster-type classifiers.

For illustrative purposes only, the learning algorithm may be considered to compute the following relationship for each pixel in the frame of the 2D image
 30 sequence

$$z_n = k_a \cdot x_n + k_b \cdot y_n + k_c \cdot R_n + k_d \cdot G_n + k_e \cdot B_n$$

where

n is the n th pixel in the key-frame image

z_n is the value of the depth assigned to the pixel at x_n, y_n

k_a to k_e are constants and are determined by the algorithm

R_n is the value of the Red component of the pixel at x_n, y_n

5 G_n is the value of the Green component of the pixel at x_n, y_n

B_n is the value of the Blue component of the pixel at x_n, y_n

This process is illustrated in Figure 1.

10 It will be appreciated by those skilled in the art that the above equation is a simplification for purposes of explanation only and would not work ideally in practice. In a practical implementation using, for example, a neural network and given the large number of pixels in an image, the network would learn one large equation containing many k values, multiplications and additions.

Apply Mapping Function to 2D Image

15 The invention next takes this mapping function and applies it to the entire frame of the 2D image sequence. For a given pixel the inputs to the mapping function are determined in a similar manner as that presented to the mapping function during the learning process. For example, if the mapping function was learnt by presenting the measurements of a single pixel as input, the mapping function will now require these same measurements for the pixels in the new image. With these inputs, the mapping function performs its learnt task and

20 outputs a depth measurement. Again, in the example for a single pixel, this depth measurement may be a simple depth value. In this example, the mapping function is applied across the entire image, to complete a full set of depth data for the image. Alternatively, if the mapping function was trained using larger sets of pixels, it is now required to generate such larger sets of pixels for the new image.

25 The higher-level measurements on these sets of pixels are made, such as mean and variance, in the same manner as that during the learning process. With these inputs now established, the mapping function produces the required depth measurement, for that set of pixels.

30 For illustrative purposes only, the algorithm determines the depth, z_n , at each pixel in the 2D image by applying the following relationship

$$z_n = k_a \cdot x_n + k_b \cdot y_n + k_c \cdot R_n + k_d \cdot G_n + k_e \cdot B_n$$

where

n is the n th pixel in the image

z_n is the value of the depth assigned to the pixel at x_n, y_n

k_a to k_e are constants previously determined by the algorithm

R_n is the value of the Red component of the pixel at x_n, y_n

G_n is the value of the Green component of the pixel at x_n, y_n

5 B_n is the value of the Blue component of the pixel at x_n, y_n

This process is illustrated in Figure 2, and results in a full depth map for the 2D image. If the resulting depth map contains regions of error, modifications may be made to the user data and the process repeated to correct these regions. The mapping function may also be applied to other frames to generate depth maps.

10 **Alternative Embodiments**

In an alternative embodiment, the user may define a set of objects and assign pixels to the objects. In this embodiment, the process of generalising the user data to the remaining pixels of the image segments the entire image into the set of objects initially identified by the user. The mapping function defining the objects or the objects themselves may be the required output of this embodiment. Alternatively, functions may be applied to the objects to specify the depth of these objects, thereby constructing a depth map for the image. These functions may take the form of depth ramps and other ways of defining the depth of objects as defined in the Applicants prior application PCT/AU00/00700.

20 In a further alternative embodiment, the training algorithm may attempt to introduce a random component into the user information. With any learning algorithm this overcomes the difficulty of over-training. Over-training refers to the situation where the learning algorithm simply remembers the training information. This is analogous to a child wrote-learning multiplication tables without gaining
25 any understanding of the concept of multiplication itself. This problem is known in the field of machine learning, and an approach to relieving the problem is to introduce random noise into the training data. A good learning algorithm will be forced to distinguish between the noise in the training data, and the quality information. In doing this, it will be encouraged to learn the nature of the data
30 rather than simply remember it. An example embodiment of this approach refers to the previous example, where the training algorithm learns the function:

$$z_n = k_a \cdot x_n + k_b \cdot y_n + k_c \cdot R_n + k_d \cdot G_n + k_e \cdot B_n$$

When presenting the inputs to the training algorithm, being x, y, R, G and B , a small noise component is added to these values. The noise component may be a small positive or negative random number.

Phase Two

5 The second phase operates on an image sequence in which some images have been identified as key-frames. It receives 3D stereo data for each key-frame, typically in the form of depth maps. The depth maps may be due to any process, such as, but not limited to, human specification, the output of the first phase described above, or some stereo information capture process for example
10 a LIDAR device. Alternatively, the 3D stereo information may be in some form other than depth maps, for example disparity information obtained from a key-frame comprising a stereo pair.

For all other frames in the 2D image sequence, the invention provides specification of the depth maps, based on the key-frame information initially
15 available. It is expected that the number of key-frames will be a small fraction of the total number of frames. Hence the invention provides a way of vastly reducing the amount of depth maps required to be initially generated.

Create Mapping Function

Once the system is provided with the key-frames and their corresponding
20 depth maps, the system then observes and analyses the key-frames and the corresponding depth map initially available, in order to create a mapping function. The mapping function may be a process or function which takes as input any given measurement of a 2D image, and provides as output a depth map for that image. This mapping is learnt by capturing the relationship between the key-
25 frame image data and depth map data available for those images.

The mapping function may take the form of any generic processing unit, where input data is received, processed, and an output given. Preferably, this processing unit is amenable to a learning process, where its nature is determined by examination of the key-frame data, and its corresponding depth map. In the
30 field of machine learning, such mapping functions are known and include, although are not limited to, neural networks, decision trees, decision graphs, model trees and nearest-neighbour classifiers.

Learn Relationships between Input Data and Desired Output Data

In a learning process, information from the 2D key-frame image is presented to the mapping function. This information may be presented on a pixel by pixel basis, where pixel measurements are provided, such as red, green and blue values, or other measurements such as luminance, chrominance, contrast and spatial measurements such as horizontal and vertical positioning in the image. Alternatively, the information may be presented in the form of higher level image features, such as larger sets of pixels and measurements on a set of pixels such as mean and variance or edges, corners etc (i.e. the response of a feature detector). Larger sets of pixels may for example represent segments in the image, being sets of connected pixels forming a homogenous region.

For illustrative purposes only, the 2D image may be represented in the form

$$x, y, R, G, B$$

where x and y represent the x and y coordinates of each pixel and R, G, B represent the red, green and blue value of that pixel.

Next, the corresponding depth map is presented to the mapping function, so that it may learn its required mapping. Normally individual pixels are presented to the mapping function, however, if higher level image features are being used, such as larger sets of pixels, or segments, the depth map may be a measurement of the depth for that set of pixels, such as mean and variance.

For illustrative purposes only, the depth map may be represented in the form

$$z, x, y$$

where x and y represent the x and y coordinates of each pixel and z represents the depth value assigned to that corresponding pixel.

The process of learning this relationship between input data, and desired output would be understood by those who have worked in the area of artificial intelligence, and may take on many forms. Preferred embodiments of a learning algorithm, are those that seek to design a mapping function which minimises some measurement of mapping error.

The learning algorithm attempts to generalise the relationships between the 2D image information and the depth map present in the key-frame examples.

This generalisation will then be applied to complete the depth maps for the entire sequence. Examples of successful learning algorithms known in the art are the back-propagation algorithm for learning neural networks, the C4.5 algorithm for learning decision trees, and the K-Means algorithm for learning cluster-type classifiers.

For illustrative purposes only, the learning algorithm may be considered to compute the following relationship for each pixel in the 2D image

$$z_n = k_a \cdot x_n + k_b \cdot y_n + k_c \cdot R_n + k_d \cdot G_n + k_e \cdot B_n$$

where

- 10 n is the n th pixel in the key-frame image
- z_n is the value of the depth assigned to the pixel at x_n, y_n
- k_a to k_e are constants and are determined by the algorithm
- R_n is the value of the Red component of the pixel at x_n, y_n
- G_n is the value of the Green component of the pixel at x_n, y_n
- 15 B_n is the value of the Blue component of the pixel at x_n, y_n

It will be appreciated by those skilled in the art that the above equation is a simplification for purposes of explanation only and would not work in practice. In a practical implementation, using for example a neural network and given the large number of pixels in an image, the network would learn one large equation

20 containing many k values, multiplications and additions.

This process is illustrated in Figure 3.

The invention next takes this mapping function and applies it across a set of 2D images that do not yet have depth maps available. For a given 2D image in that set, the inputs to the mapping function are determined in a similar manner as that presented to the mapping function during the learning process. For example,

25 if the mapping function was learnt by presenting the measurements of a single pixel as input, the mapping function will now require these same measurements for the pixels in the new image. With these inputs, the mapping function performs its learnt task and outputs a depth measurement. Again, in the example for a

30 single pixel, this depth measurement may be a simple depth value. In this example, the mapping function is applied across the entire image, to complete a full set of depth data for the image. Alternatively, if the mapping function was

trained using larger sets of pixels, it is now required to generate such larger sets of pixels for the new image. The higher-level measurements on these sets of pixels are made, such as mean and variance, in the same manner as that during the learning process. With these inputs now established, the mapping function produces the required depth measurement, for that set of pixels.

For illustrative purposes only, the algorithm determines the depth, z_n , at each pixel in the 2D image by applying the following relationship

$$z_n = k_a \cdot x_n + k_b \cdot y_n + k_c \cdot R_n + k_d \cdot G_n + k_e \cdot B_n$$

where

- 10 n is the n th pixel in the image
- z_n is the value of the depth assigned to the pixel at x_n, y_n
- k_a to k_e are constants previously determined by the algorithm
- R_n is the value of the Red component of the pixel at x_n, y_n
- G_n is the value of the Green component of the pixel at x_n, y_n
- 15 B_n is the value of the Blue component of the pixel at x_n, y_n

For a sequence of 2D images, key-frames with depth maps may be spaced throughout the sequence, in any arbitrary way. In the preferred embodiment, the mapping function will be presented with a set of key-frames, and their corresponding depth maps, which span a set of 2D images that have some commonality. In the simplest case, two key-frames are used to train the mapping function, and the mapping function then used to determine the depth maps for the 2D images between the two said key-frames. However, there is no restriction to the number of key-frames that may be used to train a mapping function. Further, there is no restriction to the number of mapping functions that are used to complete a full set of 2D images. For each pair of adjacent key-frames, a new mapping function may be used.

This process is illustrated in Figure 4.

Alternative Embodiments

In an alternative embodiment, the mapping function will be presented with a larger number of key-frames for training, with an added training variable representing the passage of time through the image sequence. Referring to the previous example, the depth value, z , for a given pixel is calculated as a function of the form:

$$z_n = k_a \cdot x_n + k_b \cdot y_n + k_c \cdot R_n + k_d \cdot G_n + k_e \cdot B_n$$

The time variable is now introduced such that this function is extended to read:

$$z_n = k_a \cdot x_n + k_b \cdot y_n + k_c \cdot R_n + k_d \cdot G_n + k_e \cdot B_n + k_f \cdot T$$

5 where:

n is the n th pixel in the image

z_n is the value of the depth assigned to the pixel at x_n, y_n

k_a to k_f are constants previously determined by the algorithm

R_n is the value of the Red component of the pixel at x_n, y_n

10 G_n is the value of the Green component of the pixel at x_n, y_n

B_n is the value of the Blue component of the pixel at x_n, y_n

T is a measurement of time, for this particular frame in the sequence

This value of time may be set to zero for the starting frame in the sequence, and set to 1.0 for the final frame in the sequence. All frames in
 15 between naturally have a time value representing their relative progress towards the final frame. When training the mapping function, the time value for the key-frames is calculated and presented to the mapping function so that the learning algorithm may incorporate this information. When the mapping function is used to convert the remaining non key-frames in the sequence, for each frame its time
 20 value is calculated, and used as input in the mapping function. This mapping function is used, as described previously, to calculate the depth information for that image.

The addition of this time variable assists the training function in generalising the information available in the key-frames. In the absence of a time
 25 variable, it is possible that the depth information in two key-frames may contradict each other. This might occur when pixels of a similar colour occur in the same spatial region in both key-frames, but belong to different objects. For example, in the first key-frame, a green car may be observed in the centre of the image, with a depth characteristic bringing it to the foreground. In the next key-frame, the car
 30 may have moved, revealing behind it a green paddock, whose depth characteristic specifies a middle ground region. The training algorithm is presented with two key-frames, that both have green pixels in the centre of the

image, but have different depth characteristics. It will not be possible to resolve this conflict, and the mapping function is not expected to perform well in such a region. With the introduction of a time variable, the training algorithm will be able to resolve the conflict by recognising that the green pixels in the centre of the image, are foreground pixels at a time near the first key-frame in the image sequence. As time progresses towards the second key-frame, the training algorithm will become more inclined to recognise green pixels in the centre of the image as the middle-ground depth of the green paddock.

The addition of the time variable also enables the algorithm to be trained with multiple key-frames. Without the time variable the key-frames may contain contradictions, due for example to new objects appearing, over a long sequence. The addition of a time variable enables the learning algorithm to perform correctly over multiple key-frames since it enables the algorithm to account for the changes in the key-frames.

For example, assuming the process is being used to calculate depth maps from a video sequence at 50 frames per second. From this sequence six key-frames, each comprising a 2D image and associated depth map, spaced 10 frames apart have been selected. The algorithm will be presented with these six key-frames, the time variable being assigned as follows:

Key-frame number	0	10	20	30	40	50
Time variable	0	0.2	0.4	0.6	0.8	1.0

These values of time would be presented to the learning algorithm for each pixel in the corresponding key-frame.

When the algorithm is interpolating the z data for frame 14 we present to the learning algorithm, for each pixel in frame 14, a time value of

$$(14/50) = 0.28$$

from this value the algorithm can determine that the frame to be calculated is in the region of key frames 10 and 20 and calculate accordingly.

In a further alternative embodiment, the training algorithm may attempt to introduce a random component into the key-frame information, as described in Phase One.

An alternative embodiment may exploit the fact that the mapping functions give a full representation of the depth information for all non key-frame images in the sequence. The mapping function could be viewed as an encoding of this depth information. It is expected that the mapping function may be transmitted with a relatively small amount of data, and hence represents a significant compression of the depth information.

Consider the case where there are two key-frames, 20 frames apart in the sequence. A mapping function has been learnt for these two key-frames, and this mapping function now provides all depth information for the intermediate frames.

The mapping function itself represents a compression of all this depth information across the twenty frames. If, for example purposes only, the mapping function can be written to a file using 6000 bytes, then for this price we gain 20 frames worth of depth information. Effectively, this represents a file size of $6000 / 20 = 300$ bytes per frame. In a practical implementation the effective compression will be substantial.

In a further embodiment, this above compression may allow for efficient transmission of 3D information, embedded in a 2D image source i.e. a 2D compatible 3D image. Since the mapping functions require a file length that is typically a tiny fraction of the 2D image data that it provides 3D information for, the addition of 3D information to the 2D image sequence is achieved with a very small overhead.

In this case, the 3D information is generated prior to viewing, or in real-time, at the viewing end, by simply applying the mapping function over each 2D image in the sequence as it is viewed. This is made possible by the fact that the types of mapping functions found in machine learning are very efficient in providing calculations *after* they have been trained. Typically the training process is slow and resource intensive, and is usually performed offline during the process of building the 3D image content. Once trained, the mapping function may be transmitted to the viewer end and will perform with a very high throughput suitable for realtime conversion of the 2D image to 3D.

The Applicant's own previous disclosures have related to techniques for converting 2D images into stereoscopic 3D images. The conversion processes

disclosed incorporated the generation of a depth map that was associated with a 2D image. In one embodiment the depth maps were created manually on a frame by frame basis. The improvement described in this application enables a fewer number of key-frames to have depth maps created and the intermediate
5 depth maps calculated. Since the key-frames represent a small fraction of the total number of frames, this new technique represents a substantial improvement in conversion efficiency both in terms of time and cost.

Other Applications

It is a specific intent of this disclosure that the invention should be applied
10 to the creation of depth maps for other than the production of stereoscopic images.

It will be known to those skilled in the art that depth maps are used extensively within the special effects industry. In order to compose live action, or computer generated images, within a 2D image it is frequently necessary to
15 manually produce a depth map for each frame of 2D image. These depth maps enable the additional images to be composed so as to appear to move with the appropriate geometry within the original 2D image.

It is also known that cameras are being developed that enable a depth map to be obtained from a live image. Typically these use laser range finding
20 techniques and are generically known as LIDAR devices. In order to capture depth maps at television frame rates an expensive and complex system is required. The application of this invention would enable simpler and less complex LIDAR devices to be constructed that need only capture depth maps at a fraction of the video field rate, or other infrequent periods, and the missing depth
25 maps produced by interpolation using the techniques described in this invention.

DATED this 29th day of September, 2000.

DYNAMIC DIGITAL DEPTH RESEARCH PTY LTD

WATERMARK PATENT & TRADEMARK ATTORNEYS
4TH FLOOR "DURACK CENTRE"
263 ADELAIDE TERRACE
PERTH WA 6000

Figure 1 – Phase One Training Process

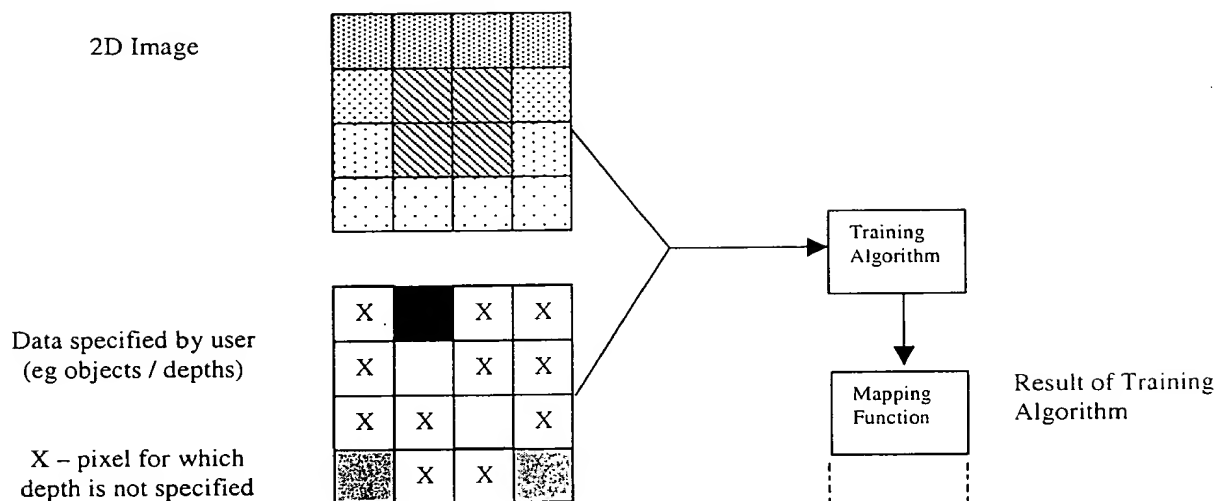


Figure 2 – Phase One Conversion Process

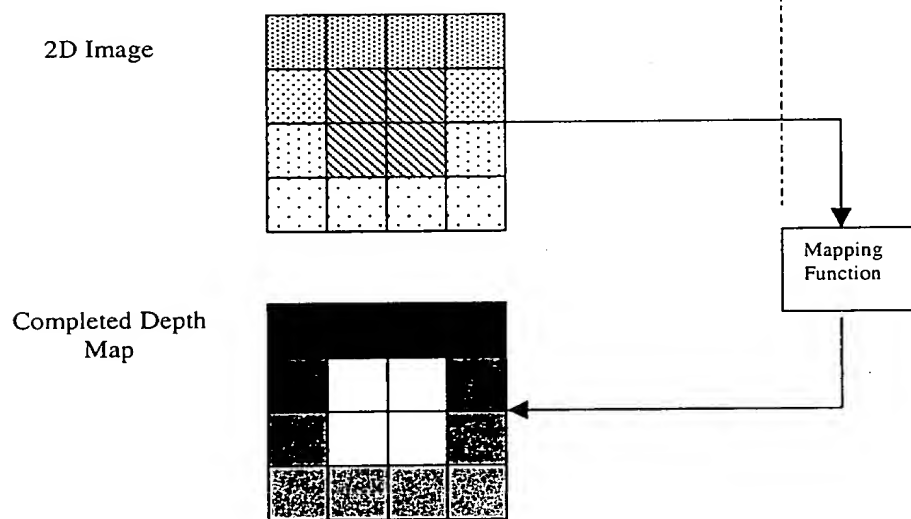




Figure 3 – Phase Two Training Process

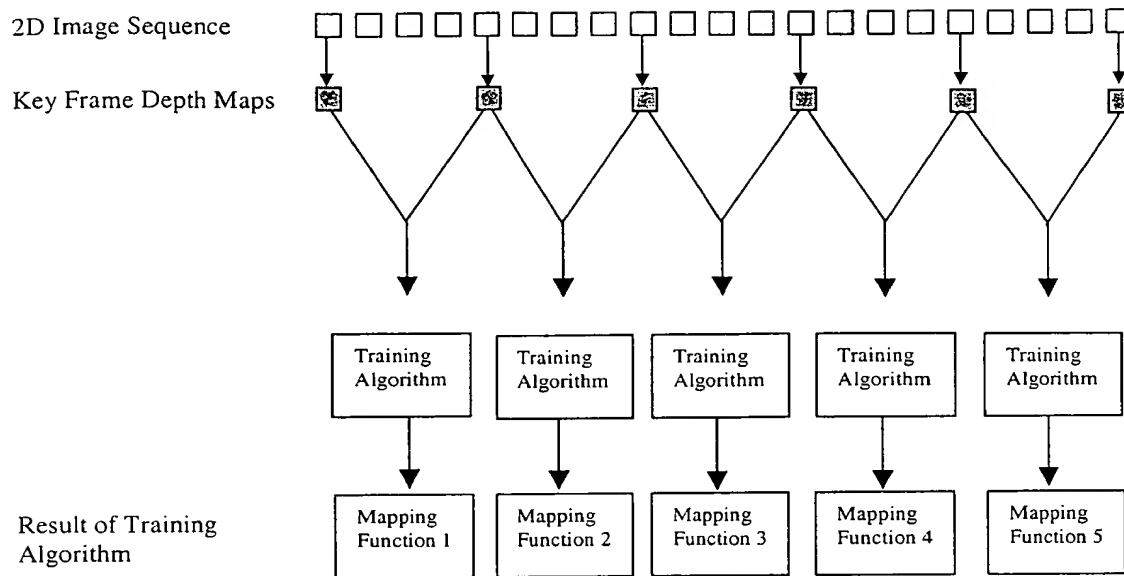


Figure 4 – Phase Two Conversion Process

